

# Piloting access to the Belgian web archive for scientific research: a methodological exploration

---

Sally Chambers, Eveline Vlassenroot, Peter Mechant and Friedel Geeraert

[Engaging with Web Archives: 'Opportunities, Challenges and Potentialities'](#) #EWAVirtual  
21-22 September 2020, Maynooth University Arts and Humanities Institute, Co. Kildare, Ireland.

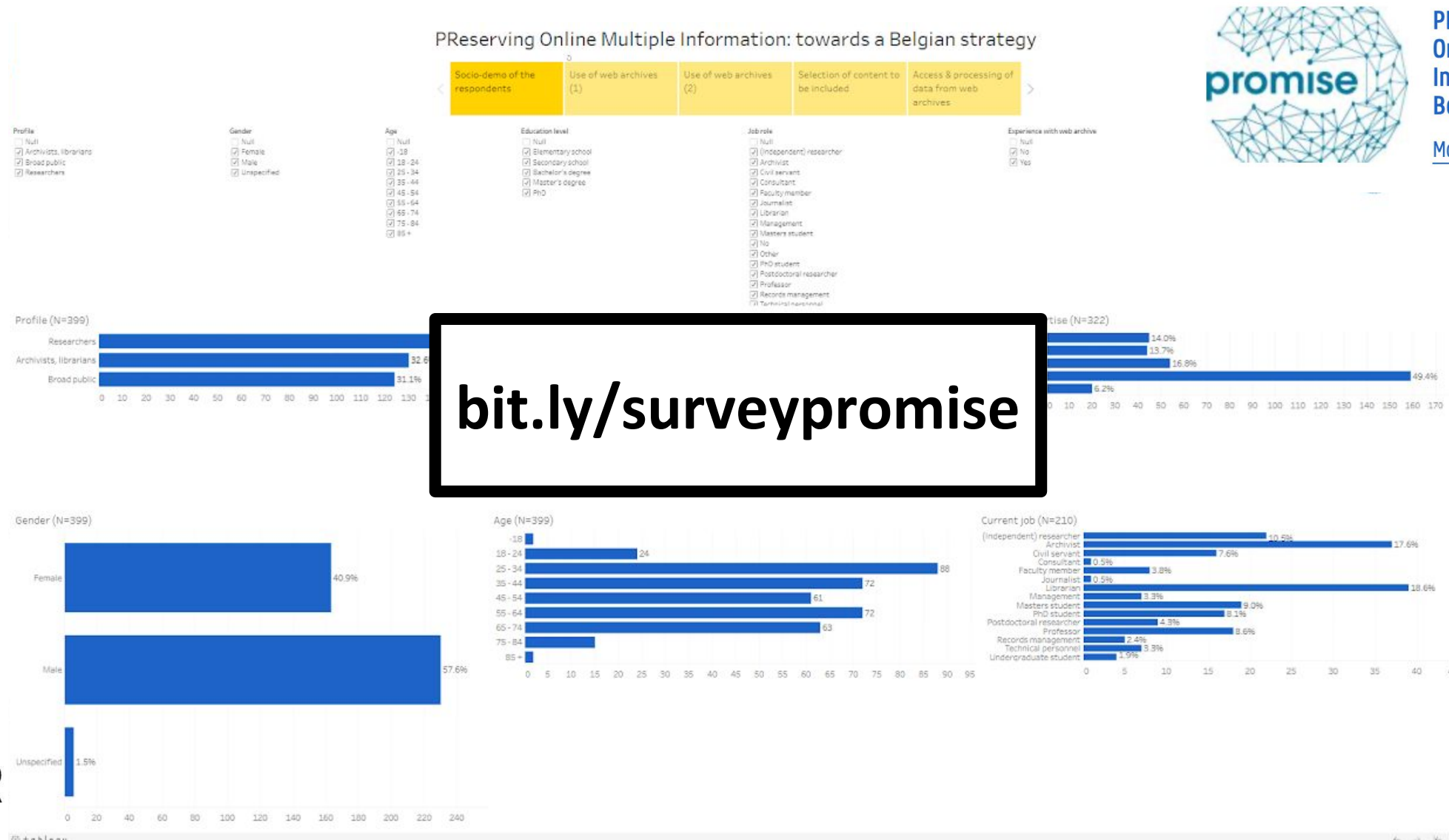
# Exploring how to pilot access to the Belgian web archive for scientific research?



- How do existing web-archives provide access to their collection for research?
- What can we learn from existing research initiatives using web archives?
- Using personas to evaluate software for the *discovery*, *delivery* & *analysis* of archived web data

[www.kbr.be/en/projects/promise-project/](http://www.kbr.be/en/projects/promise-project/)

# User requirements for a web archive



PROMISE (PReserving Online Multiple Information: towards a Belgian StratEgy)

[More information](#)

# Reviewing how existing web archives provide access to their collections for research

---

Eveline Vlassenroot – imec-MICT Ghent University

# Different types of access (Winters, 2019)

- 1) Read an archived web page
- 2) Download and manipulate data
- 3) Availability of appropriate archival frameworks and tools to facilitate navigation and research
- 4) Access facilitated by publication of research findings and data

Winters, Jane (2019). *Negotiating the archives of UK web space*. In: *The Historical Web and Digital Humanities: the Case of National Web Domains*, ed. Niels Brügger and Ditte Laursen. Routledge, Abingdon.

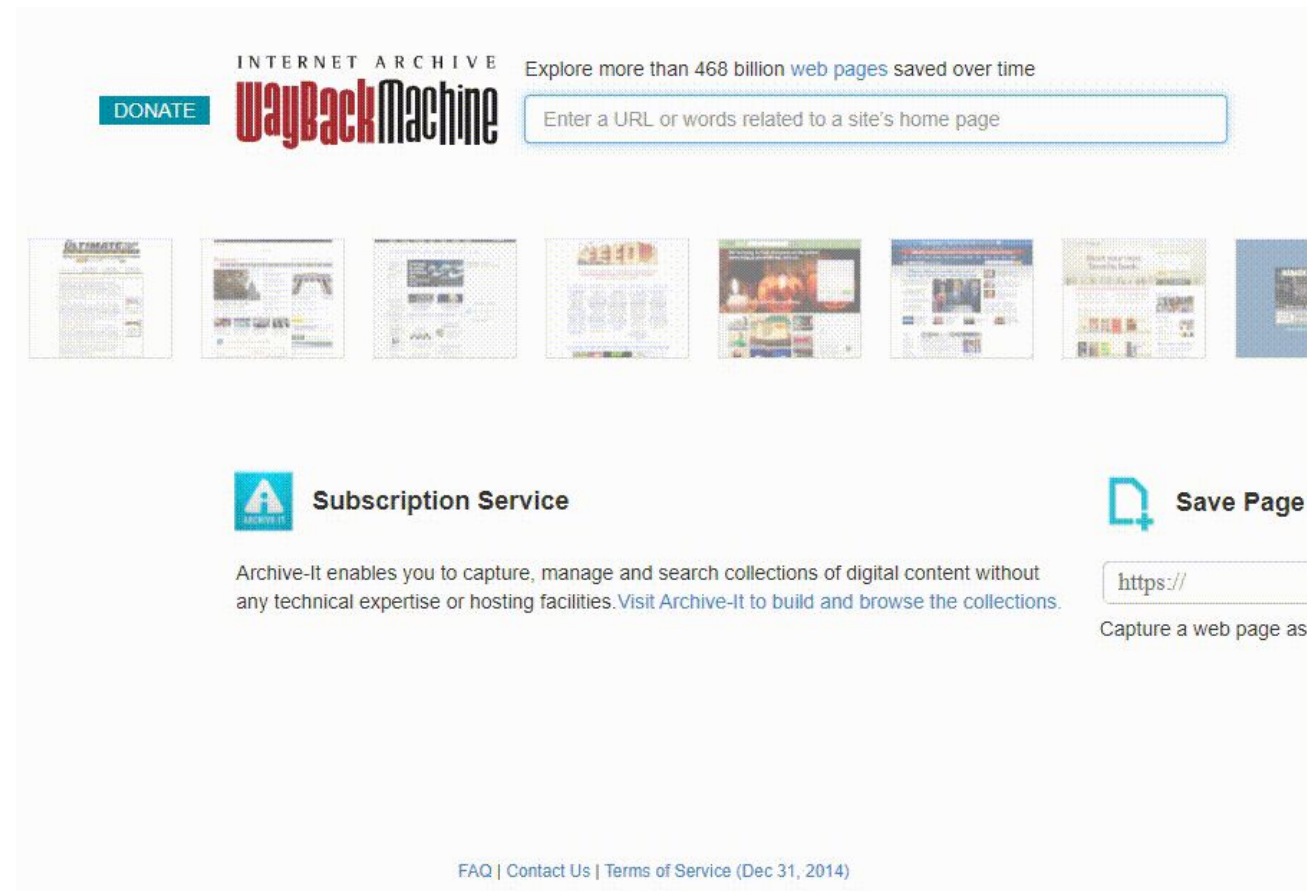
# Different types of access

- 1) Read an archived web page
- 2)
- 3)
- 4)

Country	Initiative	Web archive is available for consultation
Canada	<i>Library and Archives Canada (LAC)</i>	✓
Canada	<i>Bibliothèque et Archives nationales du Québec (BAnQ)</i>	some
Denmark	<i>Netarkivet</i>	✓
Estonia	<i>Eesti Veebiarhiiv</i>	some
France	<i>Bibliothèque nationale de France (BnF)</i>	✓
France	<i>Institut national de l'audiovisuel (INA)</i>	✓
Hungary	<i>National Széchényi Library</i>	x
Ireland	<i>National Library of Ireland</i>	✓
Luxembourg	<i>Bibliothèque nationale du Luxembourg (BnL)</i>	✓
New-Zealand	<i>National Library of New Zealand (NLNZ)</i>	✓
Switzerland	<i>Webarchiv Schweiz</i>	✓
The Netherlands	<i>KB Webarchief (KB)</i>	✓
The Netherlands	<i>Nationaal Archief (NA)</i>	some
UK	<i>UK Web Archive (UKWA)</i>	✓
USA	<i>George Washington University Libraries (GWUL)</i>	✓

# Different types of access

- 1)
- 2) Download and manipulate data
- 3)
- 4)





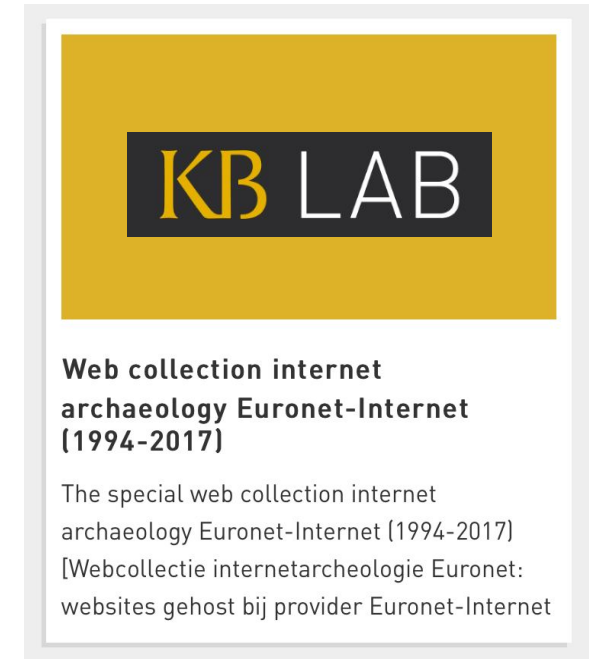
# Web archives as datasets



Web Archive Austria



Web Collection: Chinese Netherlands

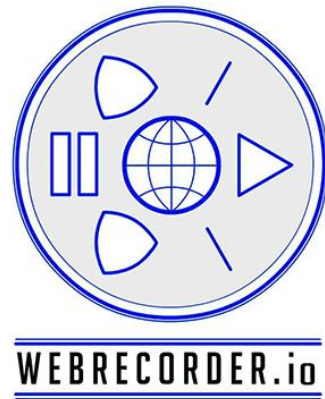


Web Collection: Internet Archaeology



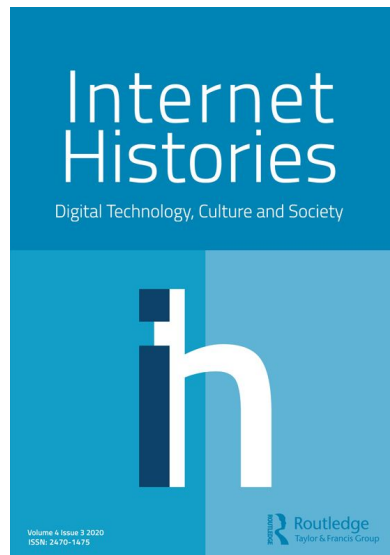
# Different types of access

- 1)
- 2)
- 3) Availability of appropriate archival frameworks and tools to facilitate navigation and research
- 4)



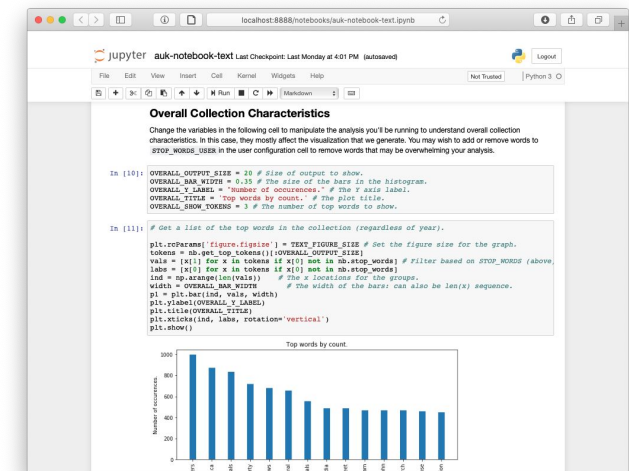
# Different types of access

- 1)
- 2)
- 3)
- 4) Access facilitated by publication of research findings and data



## Archives Unleashed Notebooks

Title	Status	Date Analyzed	Public	Files	Size
Leonard Cohen Collection			No	159	34.3 GB
University of Toronto Libraries Digital Collections			Yes	125	73.2 GB
Toronto 2015 Pan Am & Parapan American Games	Completed	April 16, 2020	Yes	294	50.4 GB
Canadian Political Parties and Political Interest Groups			Yes	6127	691 GB
Federal Election Candidate Sites 2015			Yes	310	206 GB
Toronto Mayoral Election 2014	Completed	April 16, 2020	Yes	292	292 GB
Canadian Government Information			Yes	14358	4.66 TB
Canadian Labour Unions			Yes	7757	1.03 TB
Ontario Provincial Election 2011	Completed	August 5, 2020	No	106	7.91 GB
Snowden Archive	Completed	April 16, 2020	Yes	42	716 GB
Canadian Political Interest Groups			Yes	100	8.75 GB
Ontario Provincial Election 2018	Completed	April 16, 2020	Yes	939	113 GB
University of Toronto Archives Web Collection			Yes	10624	1.35 TB
University of Toronto Scarborough			No	3	27.7 MB



## Web Archive Datasets in Zenodo & Dataverse

# What can we learn from existing research initiatives using web archives?

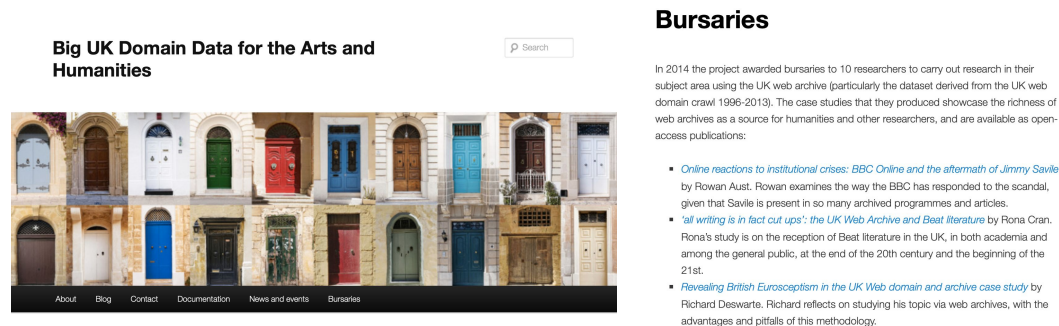
---

Sally Chambers – Ghent Centre for Digital Humanities, Ghent University

# Research-use of web archives

- **Awareness raising:** web archives as a research resource

What can we learn  
from existing  
initiatives?



- **Community building:** web archives are complex datasets: we don't need to go it alone



# Research-use of web archives

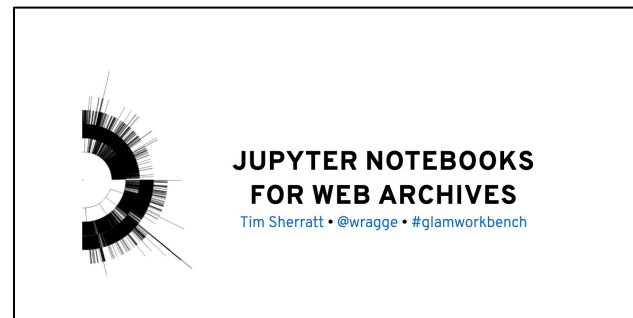
What can we learn  
from existing  
initiatives?

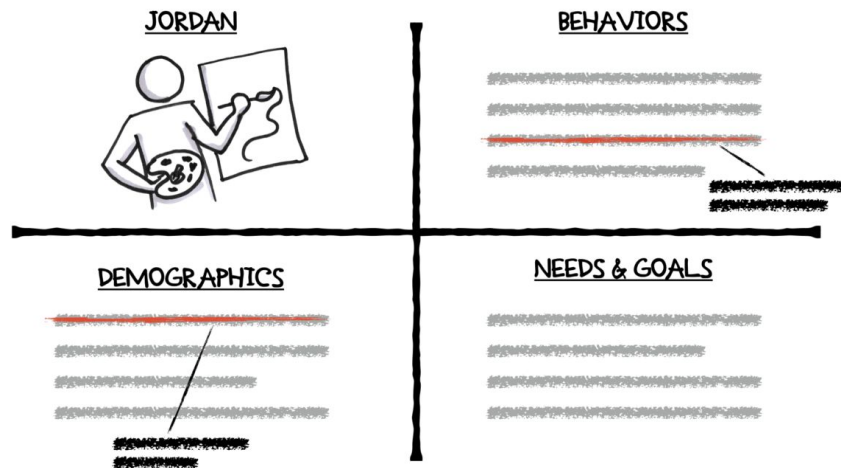
- **Research collaboration:** Research network focusing on Web Archives



**WARCnet**

- **Research infrastructure for the archived web:** sharing data, tools and workflows





Peter Mechant – imec-MICT, Ghent University

# Using **personas** to evaluate software for the Discovery, Delivery and Analysis of archived web data





Tools lifecycle -- from 2016 Harvard Environmental Scan, per Toronto bibl



Share

File Edit View Insert Format Data Tools Add-ons Help

100% \$ % .0 .00 123 Arial 10 B I S A

fx

	A	B	C	D	E	F	G	H	I	J
1		Pre-acquisition			Accessioning	Processing			Preservation	
2	Activities --->	Nomination	Rights	Assess/Define Capture	Capture	QA	Description	Indexing	Characterization	Packaging
3	Tools ↓	select sites targeted for web archiving	manage permissions to archive web sites	assess sites, define scope, create seed lists	capture web based content	enables quality assurance	add descriptive metadata	index for searching	format characterization	put into container fi
4	UNT nomination tool	1								
5	DigiBoard	1	1			1				
6	Building Collections on the Web (BCWeb)	1	1	1						
7	W3ACT	1	1	1			1			
8	Archive-It			1	1	1	1	1		
9	Web Curator Tool (WCT)		1	1	1	1	1			
10	NetarchiveSuite			1	1	1				
11	Compare Lists (of URLs)			1						
12	Extract URLs			1						
13	Expand Tiny URLs			1						
14	Harvester (list of URLs)			1						
15	<a href="#">Archiveready.com</a>			1						
16	<a href="#">Builtwith.com</a>			1						
17	Wappalzyer			1						



1 Sheet1

Sheet2



Explore

# Step 1

see Truman, Gail. 2016. Web Archiving Environmental Scan. Harvard Library Report - <http://bit.ly/1Zok3WB>  
[https://dash.harvard.edu/bitstream/handle/1/25658314/print\\_HL\\_web\\_archiving\\_env\\_scan\\_2017.pdf?sequence=4&isAllowed=y](https://dash.harvard.edu/bitstream/handle/1/25658314/print_HL_web_archiving_env_scan_2017.pdf?sequence=4&isAllowed=y)

KBR





## Le projet Corpus et ses publics potentiels.

Eleonora Moiraghi

► To cite this version:

Eleonora Moiraghi. Le projet Corpus et ses publics potentiels. : Une étude prospective sur les besoins et les attentes des futurs usagers.. [Rapport de recherche] Bibliothèque nationale de France. 2018. hal-01739730

### PERSONAS



**Alexis**

Expert de l'analyse  
de données



**Sofia**

Ingénieur de  
recherche



**Laura**

Doctorante



**Pierre**

Chercheur



**Cécile**

Professeure

see <https://hal-bnf.archives-ouvertes.fr/hal-01739730/document>

Five personas from the **CORPUS project** were considered:

## Step 2

**Cécile**, a 50-year old art history professor, with very low ICT- and software skills but with a high expertise in the domain of human and social sciences.

**Sofia**, a 34-year old research engineer, who helps researchers with all digital related issues and who has subsequently a rather high level of ICT- and software skills but rather low domain expertise in human and social sciences.

**Laura**, a 27-year old PhD-candidate in history, with mediocre ICT- and software skills and with rather high domain expertise in human and social sciences.

**Alexis**, a 28-year old data scientist, with very high ICT- and software skills but with a low expertise in the domain of human and social sciences. In a sense this persona can be contrasted with the persona of Cécile.

**Pierre**, a 36-year old researcher studying literature, who has very low ICT- and software skills but a strong expertise in the domain of human and social sciences

see <https://hal-bnf.archives-ouvertes.fr/hal-01739730/document>

# Step 3



# Step 3

	Cécile	Laura	Sofia	Pierre	Alexis
	Professor in history	PhD candidate in history	Research assistant	Researcher literature	Data scientist
Wayback Machine	+++	+++	+++	+++	+++
Archive.is	+++	+++	+++	+++	+++
Memento Time Travel	+++	+++	+++	+++	+++
Archives Unleashed Cloud	++	++	++	++	+++
Archives Unleashed Notebooks	+	++	++	+	+++
Warclight	+	++	++	+	+++
Archives Unleashed Toolkit	+	++	++	+	+++
Webrecorder pywb	+	++	++	+	+++
Browsertrix	+	++	++	+	+++
ArchiveThumbnails	+	++	++	+	+++
Shine	+	++	++	+	+++
Warcbase	+	+	++	++	+++
Natural Language Toolkit (NLTK)	+	++	+++	+	++
Leximancer	+	++	+++	+	++
WCopyfind	+	++	+++	+	+++
Mallet	+	++	+++	+	+++
CarbonDate	+++	+++	+++	+++	+++
Dorling Map Generator	+++	+++	+++	+++	+++
Raw Text to Tag Cloud Engine	+++	+++	+++	+++	+++
Gephi	+	+++	+++	+++	++

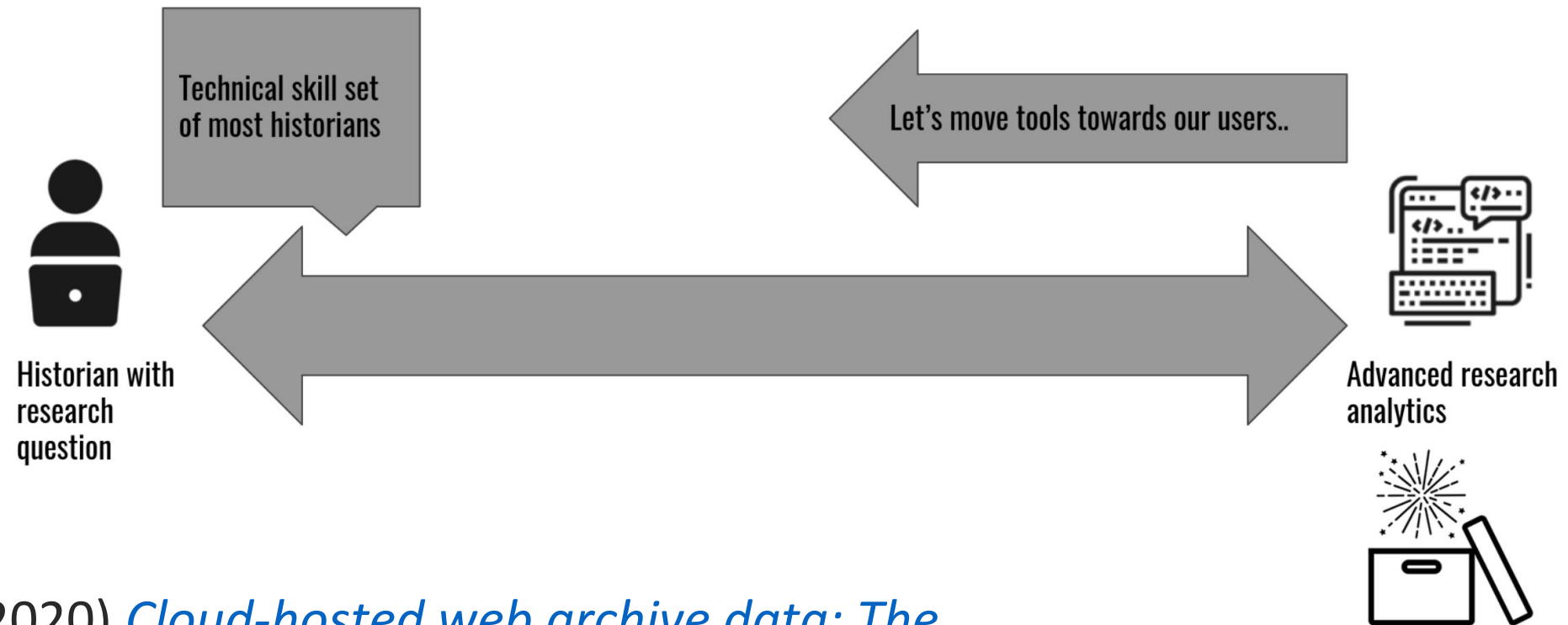
### Alexis

- as data scientist does not encounter lots of problems working with the tools as he has very high ICT- and software skills.
- due to rather low expertise in the domain of human and social sciences actually using and interpreting the results of some of these tools might pose some problems.

### Cécile

- possesses very low ICT- and software skills but has high domain expertise
- only the most basic tools offering simple (online) interfaces are in her reach;
- for more complex software she needs to acquire new skills or set up collaborations with other colleagues (e.g. Alexis or Sofia).

# Bringing web archives closer to researchers



Ruest, N. (2020) [Cloud-hosted web archive data: The winding path to web archive collections as data.](#)

Archives Unleashed project



# Exploring how to pilot access to the Belgian web archive for scientific research?



- Variety of access modes (including data-level access)
- Importance of collaboration with the web archiving community
- Increased digital literacy skills or interdisciplinary collaboration

**[www.kbr.be/en/projects/digital-research-lab/](http://www.kbr.be/en/projects/digital-research-lab/)**

# Thank you!

---

Sally Chambers, Eveline Vlassenroot, Peter Mechant and Friedel Geeraert  
(Sally.Chambers@UGent.be)

[Engaging with Web Archives: 'Opportunities, Challenges and Potentialities'](#) #EWAVirtual  
21-22 September 2020, Maynooth University Arts and Humanities Institute, Co. Kildare, Ireland.